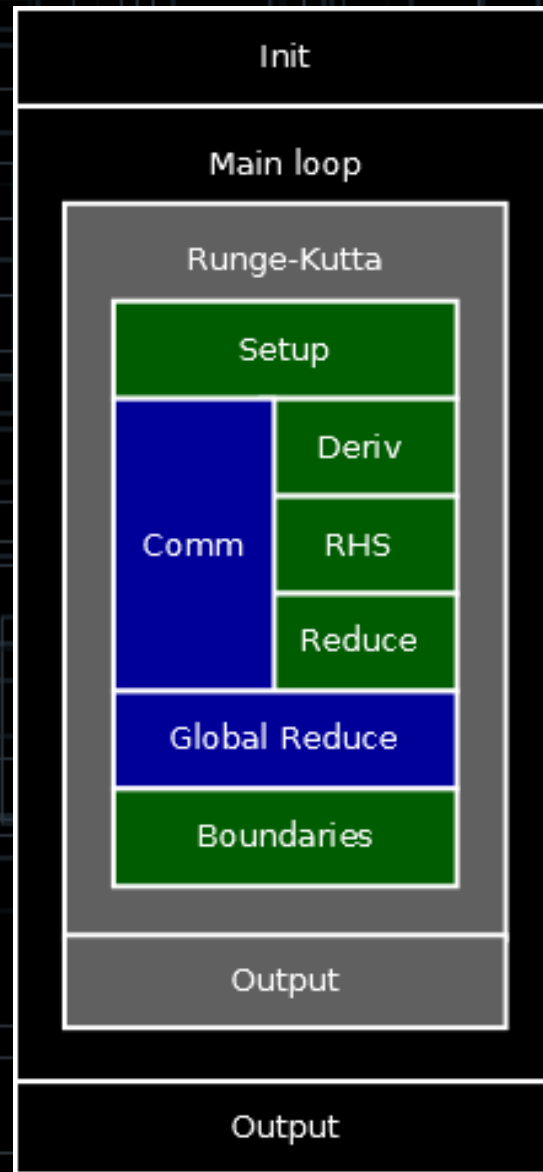


Simulación de la ecuación del flujo de Ricci en S^3 usando GPU

Carlos Bederián

14 – 11 - 2010

Estructura



X 8

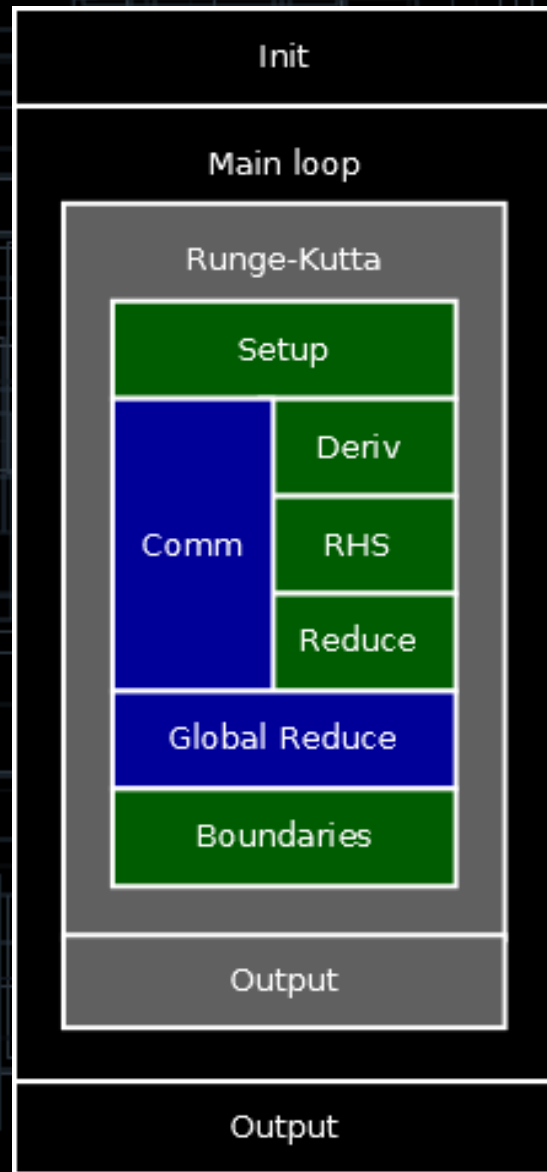
El plan

- 1.Dividir el problema en funciones
- 2.Convertir funciones en kernels
- 3.Debugging
- 4.Profiling / Mejoras
- 5.Profit

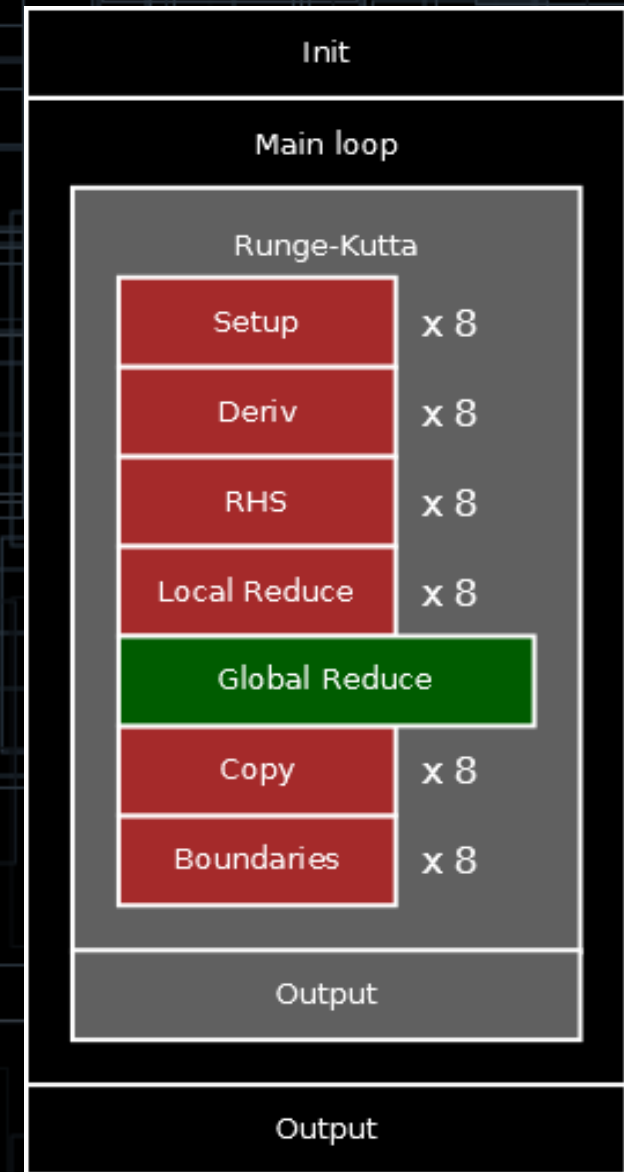
Lo que pasó

1. Dividir el problema en funciones
2. Convertir funciones en kernels
3. ~~Debugging~~
4. Eliminar MPI
5. Debugging
6. Profiling / Mejoras
7. Agregar MPI
8. Profit

Estructura



X 8



Dificultades

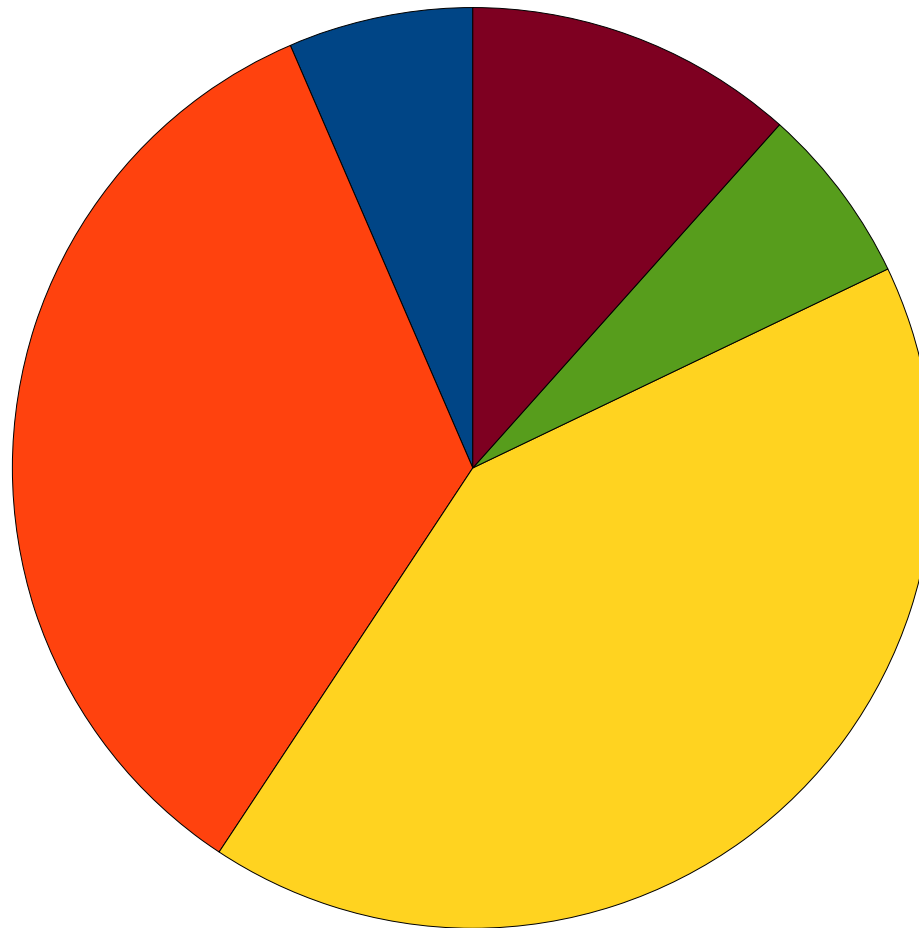
- gridDim.z debe ser 1
 - Kernels con loop sobre z
 - Correr gridDim.y x gridDim.z bloques en y
- Uso excesivo de registros
 - Occupancy baja
 - El compilador hace spilling a local memory!
 - Tratar de hacer el menor trabajo posible en estos kernels

Dificultades

- **Uso excesivo de memoria**
 - Eliminar datos duplicados
 - Recalcular campos sencillos
 - Acumular donde sea posible
 - 50% menos memoria!
- **Código multiprecisión**
 - Los literales de doble precisión eliminan toda la ventaja de usar floats
 - Regar el código con macros

¿En qué concentrarse?

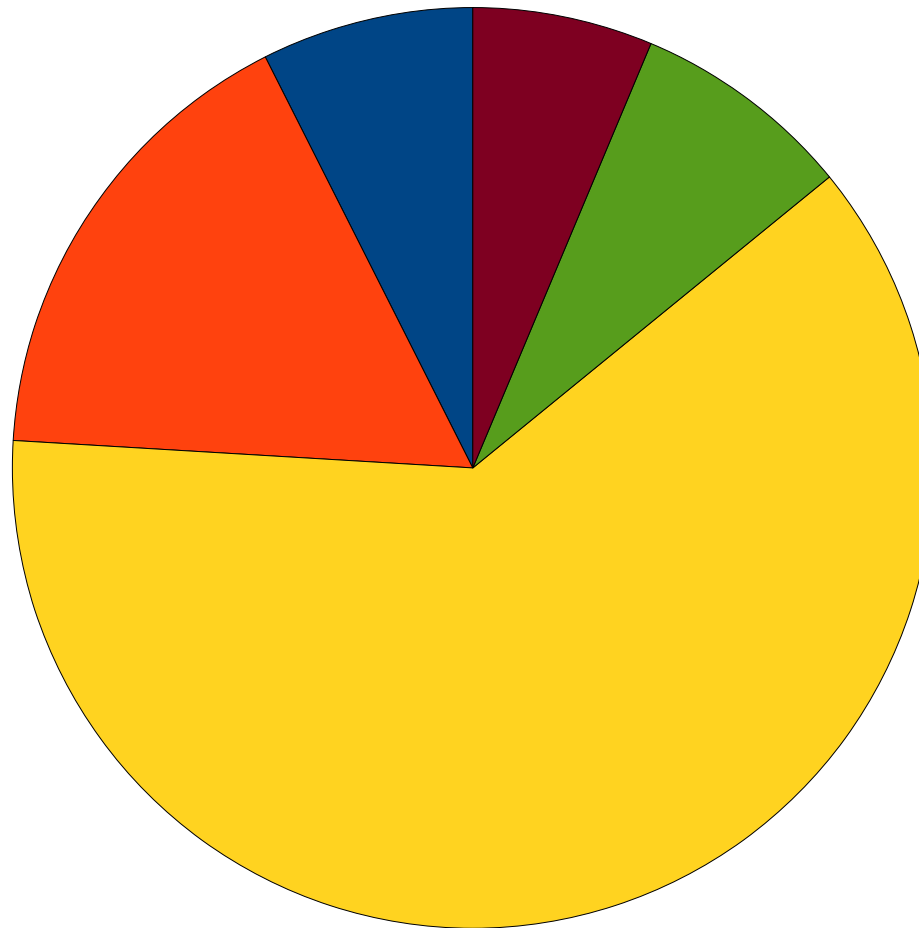
Precisión doble



■ Setup – 7% ■ RHS + Reduce – 34% ■ Deriv – 41% ■ Comm – 6% ■ Boundaries – 12%

Profile early, profile often

Precisión simple



■ Setup – 7% ■ RHS + Reduce – 17% ■ Deriv – 62% ■ Comm – 8% ■ Boundaries – 6%

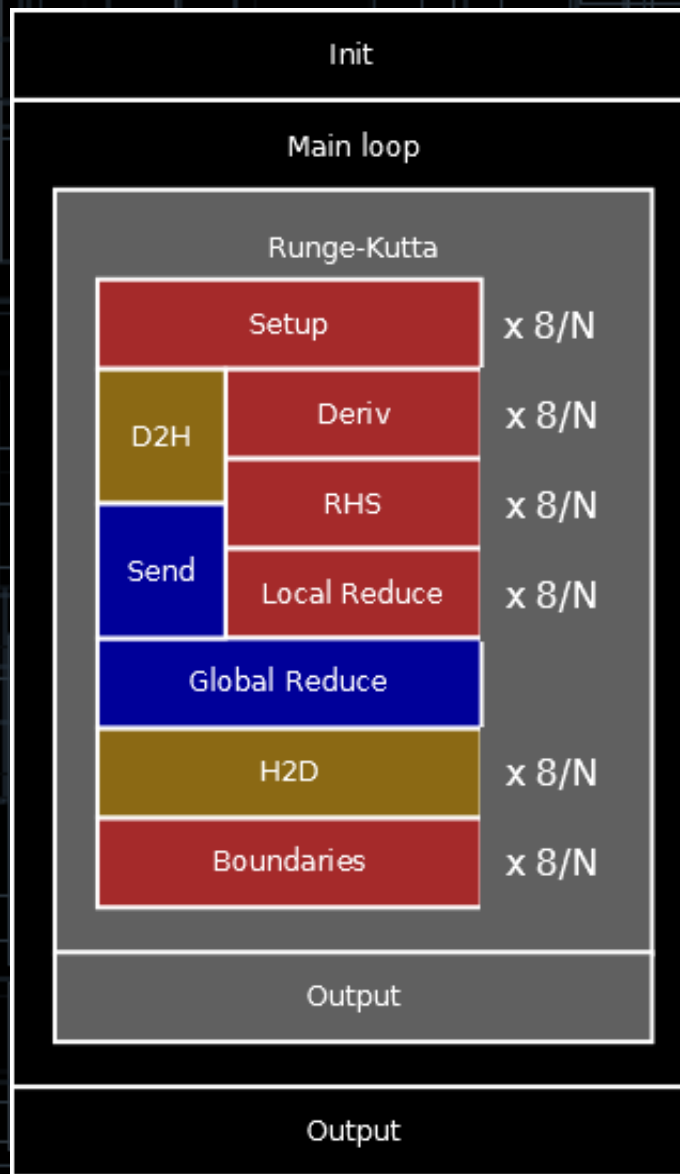
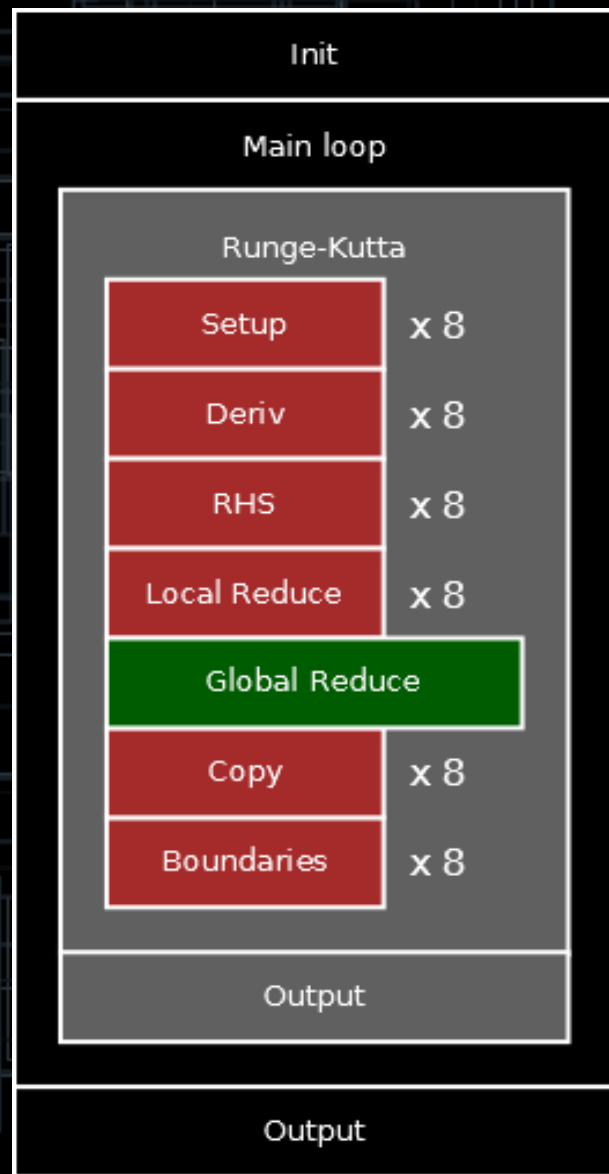
Derivadas

- 8 / 12 datos por punto
 - División entre cálculo de valores internos y externos
- Memory bound
 - Uso de shared memory
 - 25-45 GB/s en algunos casos, se puede mejorar

Reinserción de MPI

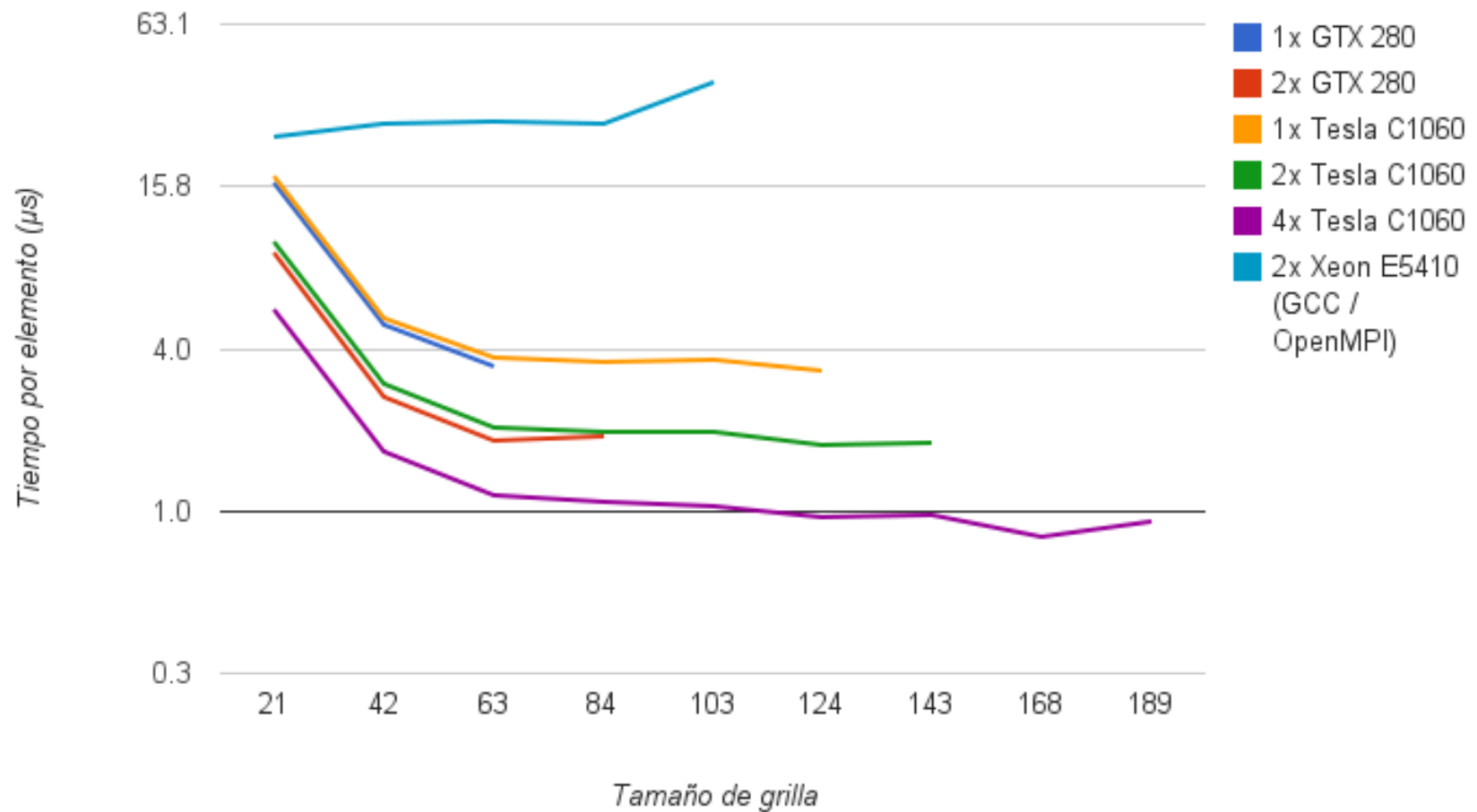
- Conservar la mayor localidad posible
 - 8 / N grids por proceso
- Código MPI aislado
 - Funciones básicas envueltas
 - Fácil de apagar para debugging
- Paralelismo de computación y comunicación
 - Uso de streams y pinned memory para copias asíncronas al host
 - Nonblocking I/O
- 2 GPUs: >1.88X, 4 GPUs: >3.44X

Estructura



X N

Resultados



Resultados

